
SpinQuant: LLM Quantization with Learned Rotations

Zechun Liu* Changsheng Zhao* Igor Fedorov Bilge Soran Dhruv Choudhary
Raghuraman Krishnamoorthi Vikas Chandra Yuandong Tian Tijmen Blankevoort
Meta

Abstract

Post-training quantization (PTQ) techniques applied to weights, activations, and the KV cache greatly reduce memory usage, latency, and power consumption of Large Language Models (LLMs), but may lead to large quantization errors when outliers are present. Rotating activation or weight matrices helps remove outliers and benefits quantization. In this work, we identify a collection of applicable rotation parameterizations that lead to identical outputs in full-precision Transformer architectures while enhancing quantization accuracy. In addition, we find that some random rotations lead to much better quantization than others, with an up to *13 points* difference in downstream zero-shot reasoning performance. As a result, we propose SpinQuant, a novel approach that incorporates learned rotation matrices for optimal quantized network accuracy. With 4-bit quantization of weight, activation, and KV-cache, SpinQuant narrows the accuracy gap on zero-shot reasoning tasks with full precision to merely 2.9 points on the LLaMA-2 7B model, surpassing LLM-QAT by 19.1 points and SmoothQuant by 25.0 points. Furthermore, SpinQuant also outperforms concurrent work QuaRot, which applies random rotations to remove outliers. In particular, for LLaMA-3 8B models that are hard to quantize, SpinQuant reduces the gap to full precision by up to 45.1% relative to QuaRot.

1 Introduction

Large Language models (LLMs) have demonstrated impressive performance across many disciplines. SoTA open source models (*e.g.*, LLaMA [41], Mistral [17], etc) and proprietary LLMs (*e.g.*, GPT [2], Gemini[38], etc) have been used in general purpose chatting assistants, medical diagnosticians [39], computer game content generators [10], coding co-pilots [34], and much more.

To serve such a high demand, the inference cost becomes a real issue. Many effective techniques have been developed. Post-training Quantization (PTQ), as one effective category of techniques, quantizes the weights (or activations) into low-precision and thus reduces the memory usage and may significantly improve latency. This is not only important for server-side inference, but also for on-device scenarios with small-sized LLMs [27, 3].

When applying quantization, **outliers remain an open challenge because they stretch the quantization range, leaving fewer effective bits available for the majority of values.** Prior research mitigates this challenge by trading quantization difficulty between weights and activations [46, 23] or employing mixed-precision to handle outliers [50]. In this work, **we focus on a new angle: multiplying the weight matrix with a rotation matrix to reduce outliers and enhance quantizability.** Inspired by [13] and SliceGPT [4], we leverage the property of rotational invariance to **construct rotation matrices in pairs from identity mapping, which can be integrated into nearby weights without affecting the overall network outputs.** By applying these random rotations, we produce a distribution of weight or activation entries that is outlier-less, facilitating easy quantization.

* Equal contribution. Correspondence to: Zechun Liu <zechunliu@meta.com>.

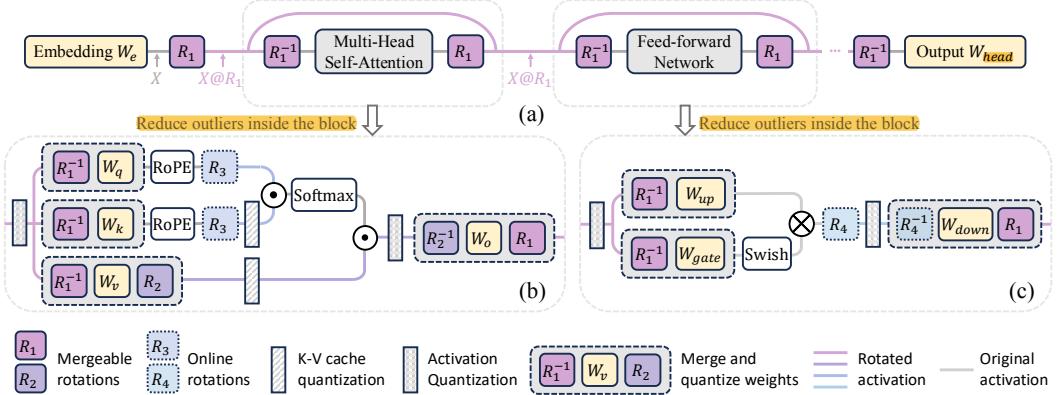


Figure 1: **Overall diagram of rotation.** (a) The residual stream can be rotated in the transformer network, resulting in numerically equivalent floating point networks before and after rotation. The rotated activations exhibit fewer outliers and are easier to quantize. (b) & (c) The rotation matrix can be integrated with the corresponding weight matrices and we further define R_2 , R_3 , and R_4 for reducing outliers inside the block.

In addition to using random rotation, which statistically works well, we find that the performance of quantized network could *vary a lot* with different rotation matrices. For example, the downstream averaged accuracy on zero-shot reasoning tasks may change up to 13 points with different rotations. As a result, we propose SpinQuant that *integrates* and *optimizes* the rotation matrix to minimize the final loss of the quantized network, with fixed weight parameters, by employing the *Cayley SGD* [21], a proficient technique for optimizing orthonormal matrices. This optimization does not alter the full-precision network output but refines the intermediate activations and weights, making them more quantization-friendly.

In SpinQuant, we introduce two rotation strategies tailored for different complexity levels: SpinQuant_{no had} and SpinQuant_{had}. Here, *had* refers to hadamard rotation matrix. In SpinQuant_{no had}, as depicted in Figure 1(b), we implement shortcut rotation (R_1) and W_v - W_o pair rotation (R_2), which can be directly absorbed into the respective weight matrices. During inference, the original weights are simply replaced with the rotated quantized weights, eliminating the need for modification in the forward pass. Conversely, in SpinQuant_{had}, designed for scenarios with low-bit quantization of KV cache or activations (e.g., 4-bit), we further incorporate online Hadamard rotation matrices (R_3, R_4) to address activation outliers inside MLP block and KV cache.

To rigorously assess the effectiveness of SpinQuant, we executed comprehensive experiments across seven leading Large Language Models (LLMs), including LLaMA-2[41] models (7B/13B/70B), LLaMA-3[3] models (1B/3B/8B), and the Mistral [17] 7B model. The key contributions of this study are summarized as follows:

- We introduce SpinQuant, the first method that employs learned rotations to mitigate outliers in weight and activation distributions, boosting the performance of quantized LLMs.
- We reveal that random rotations introduce substantial variance in quantized network performance. We propose optimizing rotation matrices within *Stiefel* manifold, directly minimizing the final loss of rotated quantized network. Ablation studies validate that our learned rotations consistently outperform random rotations, with improvements up to 16.2 points.
- SpinQuant_{no had} merges rotation matrices into pre-trained weights without altering the network architecture, significantly narrowing the W4A8KV8 quantization performance gap from 12.1 to 1.6 on the Mistral-7B model in zero-shot commonsense reasoning tasks. Noteworthily, SpinQuant_{no had} W4A8 quantization achieves comparable performance as state-of-the-art weight-only quantization methods like QuIP# [42] and OminiQuant [37] on LLaMA-2.
- SpinQuant_{had} attains an average accuracy of 64.0 in extreme W4A4KV4 quantization settings on LLaMA-2 7B. This represents a mere 2.9 point gap from the full-precision network, a substantial improvement over the previous LLM-QAT [26] approach, which exhibited a 22.0 point gap under identical precision conditions.

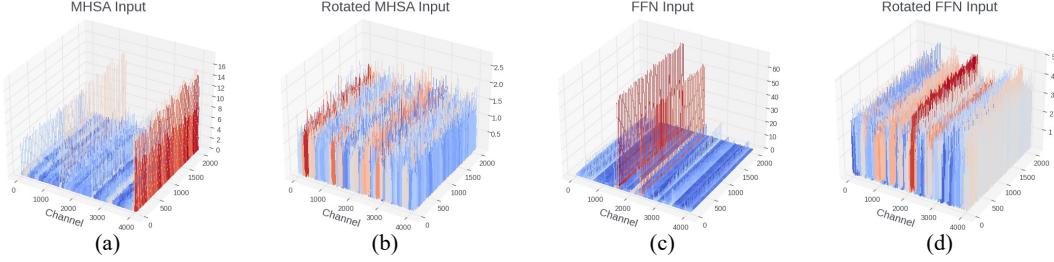


Figure 2: Activation distribution in LLaMA-2 7B model before and after rotation. Outliers exist in particular channels before rotation. Since channel-wise quantization is not supported in most hardware, outlier removal using rotation enables accurate token-wise or tensor-wise quantization.

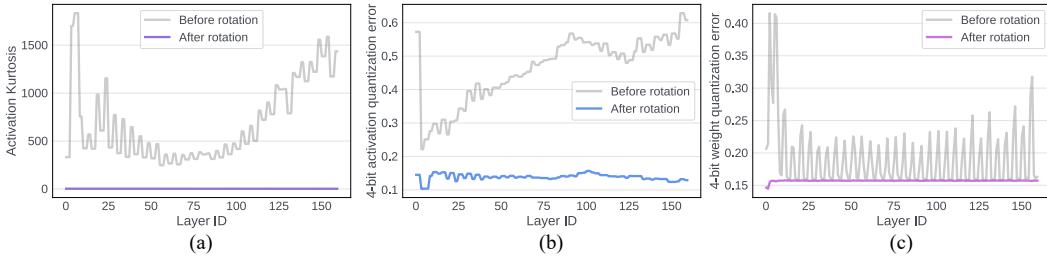


Figure 3: Outlier measurement and quantization error across input activation and weights in the five layers that take inputs from the residual (Q/K/V/Up/Gate-projection) of each block in the LLaMA-2 7B model. (a) After rotation, *kurtosis* of activation distributions is significantly reduced to approximately three across all layers. Quantization error is reduced after rotation in both (b) activations and (c) weights.

2 Motivation

Quantization reduces the precision of weights (and/or activations) in a neural network in order to save memory and lower the latency. The quantization process can be formulated as:

$$X_Q = \alpha \lfloor \frac{X_R - \beta}{\alpha} \rfloor + \beta \quad (1)$$

where $\alpha = \frac{\max(|X_R|)}{2^{N-1}-1}$, $\beta = 0$ in symmetric quantization or $\alpha = \frac{\max(X_R)-\min(X_R)}{2^N-1}$, $\beta = \min(X_R)$ in asymmetric quantization. Here X_Q is a quantized tensor and X_R is a real-valued FP16 tensor. N is number of bits. For Large language models (LLMs), the presence of outliers extends the range of weight/activation values and increases the reconstruction errors for normal values [11, 25, 48] (Figures 2 (a)&(c)).

2.1 Outlier Reduction

There exist many ways to mitigate the effect of outliers [46, 11]. In this paper, we propose to use optimized rotation to reduce outliers. Intuitively, a random rotation matrix statistically blends large and small weights together into a well-behaved distribution with fewer outliers [13], and thus is easier to quantize.

Figure 3 (a) illustrates the measurement of the *Kurtosis* κ of the activations before and after rotation. κ quantifies the “tailedness” of a real-valued random variable’s probability distribution. A larger κ indicates more outliers, while $\kappa \approx 3$ suggests a Gaussian-like distribution. In Figure 3 (a), the activation distribution in the transformer contains numerous outliers, with κ of many layers exceeding 200. However, after multiplying these activations with a random rotation matrix, the κ across all layers becomes approximately 3, indicating a more Gaussian-shaped distribution that is easier to quantize. This is corroborated by Figure 3 (b), where the quantization error of the activation tensor significantly decreases after rotation.

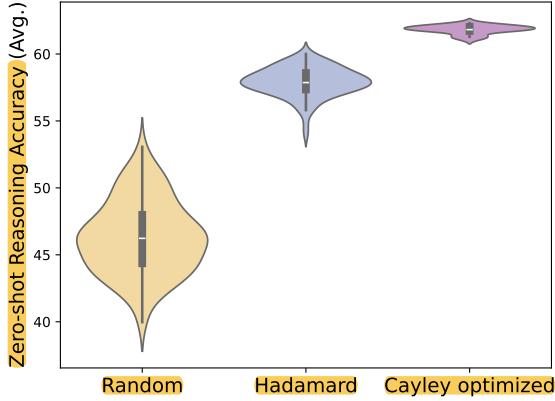


Figure 4: The performance distributions of W4A4 quantized LLaMA-2 7B under different random rotations, using network-level parameterization (Sec. 3.1). We compare the distributions using random floating-point rotations, random Hadamard matrices, and optimized rotation matrices with *Cayley* optimization (Sec. 3.2). Despite that Hadamard matrices mostly perform better than random rotations, both random groups demonstrate large variance. In contrast, by optimizing the rotation matrix with *Cayley* optimization (*i.e.*, SpinQuant), the performance is improved significantly and the variance becomes much smaller.

2.2 Random rotations produce large variance

Interestingly, while statistically random rotation leads to better quantization, not all random rotations give the same quantization outcome. To show this, we tested the zero-shot average accuracy of the rotated version of LLaMA-2 7B, quantized to 4-bit weight and 4-bit activation, under 100 randomized trials. As shown in Figure 4, the performance variance is substantial, with the best random rotation matrix outperforming the worst by 13 points. Random Hadamard matrices² outperform random rotation matrices, in consistent with the findings in [42] that Hadamard matrices yield tighter bounds on weight maximal value. However, even random Hadamard rotation matrices exhibit a non-negligible variance in final performance, as large as 6 points.

Given the huge variance across multiple trials of rotations, a natural question arises: *Is it possible to optimize the rotation to maximize the benefit of quantization?* We affirmatively answer this question by presenting a viable framework with quantization-oriented rotation learning that consistently achieves high accuracy across 7 models and 4 low-bit quantization settings.

3 Method

In this section, we introduce SpinQuant, a framework that integrates and optimizes rotations in LLMs targeting at quantization loss. We start with defining rotation parameterization of popular LLM architectures, which includes two mergeable rotation matrices (R_1, R_2) that produce rotationally invariant full-precision network, and two online Hadamard rotation (R_3, R_4) to further reduce the outliers for extreme activation and KV-cache quantization. Then, we present how to optimize these rotation matrices on *Stiefel* manifold with target loss.

3.1 Rotation parameterization

Rotating activations in residual As shown in Figure 1(a), we rotate the activations in the residual path by multiplying the embedding output X with a random rotation matrix (R_1). This rotation removes outliers and eases the quantization of the input activations to the fully-connected layers that read from the residual. To maintain numerical invariance, we reverse the rotation of the activation by

²A Hadamard matrix H is a special type of rotation matrix, where the entries of the matrix are solely $\pm\sqrt{n}$. Given a Hadamard matrix H , we can generate 2^n different random Hadamard matrices by multiplying with S , a diagonal matrix with elements s_i randomly chosen from $\{-1, 1\}$.

multiplying it with R_1^T ($= R_1^{-1}$) prior to its passage through the attention block and feed-forward network, which contains non-linearity. When the quantization is not present, the full-precision network remains intact no matter which rotation is applied.³ The rotation matrices can be merged into corresponding weight matrices, as illustrated in Figures 1(b)&(c). After absorption, no new parameters are introduced in the network. We can now modify R_1 freely without impacting the floating-point network’s accuracy or parameter count.

Rotating activations in the attention block As depicted in Figure 1(b), in the attention block, we propose to rotate the value matrix by multiplying R_2 , and the activations to out-projection layer by R_2^T head-wisely. R_2 has the shape of $(D_{\text{head}}, D_{\text{head}})$ and can be independently chosen across layers. The numerical in-variance is illustrated in Figure 5, these two rotations can be offset in a full-precision network since there are no operators between R_2 and R_2^T . Meanwhile, it can improve quantization for value cache and input activations to out-projection layer without introducing any new parameters in the network.

$$\text{out} = \text{Attn} \cdot X \cdot \boxed{\boxed{W_v} \boxed{R_2} \boxed{W_o}}$$

Figure 5: Rotation equivalence in Multi-Head Self-Attention.

We denote the method with only R_1 and R_2 inserted and optimized as $\text{SpinQuant}_{\text{no had}}$, which can readily achieve significant accuracy improvement than previous quantization methods, and closing the gap between W4A8 quantized LLMs and their full-precision counterparts to $0.1 - 2.5$ points on zero-shot commonsense reasoning averaged accuracy.

Additional unabsorbed rotations To further enhance outlier suppression for lower-bit (*e.g.* 4-bit) activation quantization, we incorporate a Hadamard matrix multiplication (R_4 in Figure 1(c)) inside the feed-forward block, reducing the outliers in the input to the down projection layer, similar to [42, 5]. Hadamard rotation can be computed with fast hadamard transform and introduce marginal overhead to the inference latency. Similarly, Hadamard matrix (R_3 in Figure 1(b)) can be inserted when low-bit KV cache quantization is required. We denote the resulting method, equipped with all rotations, as $\text{SpinQuant}_{\text{had}}$. Next, we demonstrate how to jointly optimize these rotations.

3.2 Cayley-optimized rotation

As illustrated in Figure 1, we have determined that the incorporation of four rotation matrices (R_1, R_2, R_3, R_4) can improve quantization performance while preserving numerical consistency in a full-precision network. Given that R_3 and R_4 are online rotation operations, meaning they cannot be absorbed into the weight matrix, we retain them as Hadamard matrices. This is because online Hadamard transforms can be efficiently implemented without significant overhead. We then define the optimization objective as identifying the optimal rotation matrix R_1 and R_2 that minimizes the final loss of the quantized network:

$$\arg \min_{R \in \mathcal{M}} \mathcal{L}_Q(R_1, R_2 | W, X) \quad (2)$$

Here, \mathcal{M} represents the *Stiefel* manifold *i.e.*, the set of all orthonormal matrices. $\mathcal{L}_Q(\cdot)$ denotes the task loss, such as cross-entropy, on the calibration set. It is a function of $\{R_1, R_2\}$, given the fixed pretrained weights W and the input tensor X and with the quantization function Q in the network. To optimize the rotation matrix on the *Stiefel* manifold, we employ the *Cayley SGD* method [21], which is an efficient optimization algorithm on the *Stiefel* manifold. More specifically, in each iteration, the update of the rotation R is parameterized as the following:

$$R' = \Delta R(Y)R := \left(I - \frac{\alpha}{2}Y\right)^{-1} \left(I + \frac{\alpha}{2}Y\right)R \quad (3)$$

³In a pre-norm LLM like LLaMA [40], we can convert a transformer network into a rotation-invariant network by incorporating the RMSNorm scale parameters α into the weight matrix right after the RMSNorm layer [4].

where $\Delta R(Y) := (I - \frac{\alpha}{2}Y)^{-1}(I + \frac{\alpha}{2}Y)$ is the *Cayley Transform* of a skew-symmetric matrix Y (i.e., $Y^\top = -Y$). Y is computed from a projection \hat{G} of the gradient $G := \nabla_R \mathcal{L}_Q$ of the loss function:

$$Y = \hat{G} - \hat{G}^\top, \quad \hat{G} := GR^\top - \frac{1}{2}RR^\top GR^\top \quad (4)$$

It can be shown that $\Delta R(Y)$ is always orthonormal and thus R' is guaranteed to be orthonormal ($R'^\top R' = I$) if R is orthonormal. While Eqn. 3 requires a matrix inverse, the new rotation matrix R' can be computed via an efficient fixed point iteration [21]. Overall, the approach maintains the property of orthonormality with only ~ 2 times the computation time per iteration compared to a naive SGD algorithm.

We apply the *Cayley SGD* method to solve Eqn. 2 for $\{R_1, R_2\}$, while the underlying weight parameters in the network remain frozen. $\{R_1, R_2\}$ count for only $\sim 0.26\%$ of the weight size and is constrained to be orthonormal. Consequently, the underlying floating-point network remains unchanged, and the rotation only influences the quantization performance.

By employing *Cayley* optimization to update the rotation for 100 iterations on an 800-sample WikiText2 calibration dataset, we obtain a rotation matrix that outperforms the best random matrix and random Hadamard matrix in 100 random seeds, shown in Figure 4. The *Cayley*-optimized rotation exhibits minimal variance when initiated from different random seeds. The rotation matrices are initialized with random Hadamard matrices for optimization and our ablation study in Section 4.3.3 demonstrates that the optimized rotation is robust to random rotation initialization as well.

4 Experiments

We conduct experiments on the LLaMA-2 [41] models (7B/13B/70B), LLaMA-3 [3] models (1B/3B/8B) and Mistral [17] 7B model. Our evaluation of the proposed SpinQuant was carried out on eight zero-shot commonsense reasoning tasks. These tasks include BoolQ [8], PIQA [6], SIQA [36], HellaSwag [49], WinoGrande [35], ARC-easy and ARC-challenge [9], and OBQA [29]. Additionally, we also report the perplexity score on WikiText2 testset [28] for our evaluation.

4.1 Experimental settings

We employ *Cayley SGD* [21] to optimize the rotation matrix, R_1 and R_2 , both initialized as a random Hadamard matrix, while maintaining all network weights constant. R_1 is the residual rotation, shaped as $(D_{\text{token}}, D_{\text{token}})$. R_2 is head-wise rotation in each attention block, shaped as $(D_{\text{head}}, D_{\text{head}})$ and is separately learned in each layer. The learning rate starts at 1.5 and linearly decays to 0. We utilize 800 samples from WikiText-2 to optimize rotation for 100 iterations. It takes only $\sim 13 / 18 / 30$ minutes for LLaMA-3 1B / 3B / 8B, respectively, and $\sim 25 / 30$ minutes for LLaMA-2 7B / 13B, respectively. For LLaMA-2 70B, it takes ~ 3.5 hours and for Mistral-7B it takes ~ 16 minutes.

In the main results, we optimize the rotation with respect to the activation quantized network, where the weights remain 16-bit. After rotation is learned, we apply GPTQ on the rotated weights [14], for which we adhere to the standard GPTQ settings by using 128 samples from WikiText-2 with a sequence length of 2048 as the calibration set for GPTQ quantization. In the main table, we present the results of SpinQuant with GPTQ, and in the ablation study, while we also show the results of employing simple round-to-nearest (RTN) quantization in the ablation study.

4.2 Main results

We present two rotation schemes $\text{SpinQuant}_{\text{no had}}$ and $\text{SpinQuant}_{\text{had}}$ to accommodate different scenarios. In Table 1, we use seven models and four most commonly used bit-width settings to provide a guideline on which rotation scheme should be chosen in practice.

Recap $\text{SpinQuant}_{\text{no had}}$ uses learned rotation R_1 and R_2 only, which can be merged into corresponding model weights during inference time after the rotation is learned. Using $\text{SpinQuant}_{\text{no had}}$ only needs to replace the original model weights with the rotated model weights, necessitating no modification to the forward pass nor any additional kernel support. While $\text{SpinQuant}_{\text{had}}$ comprises both learned rotations (R_1, R_2) and the online Hadamard rotations (R_3, R_4). During inference time,

Table 2: Compared to Hadamard rotation, SpinQuant learned rotation consistently outperform by a significant margin. Results are averaged accuracy on eight Zero-shot CommonSense Reasoning tasks.

	LLaMA-3.2 3B		LLaMA-3 8B		Mistral-7B	
	4-4-16	4-4-4	4-4-16	4-4-4	4-4-16	4-4-4
Random Hadamard $R_{\{1,2\}}$	49.8	49.6	49.5	50.0	51.4	51.5
SpinQuant _{no had} $R_{\{1,2\}}$	52.9 ($\uparrow 3.1$)	52.9 ($\uparrow 3.3$)	51.9 ($\uparrow 2.4$)	52.6 ($\uparrow 2.5$)	52.7 ($\uparrow 1.3$)	52.4 ($\uparrow 0.9$)
Random Hadamard $\bar{R}_{\{1,2,3,4\}}$	-	-	59.0	58.4	64.2	63.9
SpinQuant _{had} $R_{\{1,2,3,4\}}$	61.0 ($\uparrow 2.1$)	60.5 ($\uparrow 2.2$)	65.8 ($\uparrow 1.6$)	65.5 ($\uparrow 1.6$)	68.4 ($\uparrow 15.7$)	68.6 ($\uparrow 16.2$)

Table 3: Ablation study on compatibility with GPTQ [14] on a LLaMA2-7B model.

#Bits _(W-A-KV)	Task	Cayley on 4-4-KV	Cayley on 16-4-KV
4-4-16	0-shot ⁸ Avg. Wiki	61.0 ± 1.0 6.7 ± 0.07	64.1 ± 0.4 5.9 ± 0.00
4-4-4	0-shot ⁸ Avg. Wiki	60.9 ± 0.6 6.8 ± 0.15	64.0 ± 0.3 5.9 ± 0.01

accuracy across various models and bit-width configurations. Notably, in the quantization of Mistral-7B, SpinQuant_{had} secures an improvement exceeding 15.7 points over using random Hadamard rotations. Given that rotation optimization incurs a minimal time cost (only 30 minutes for smaller models and up to 3.5 hours for a 70B model) we advocate for the adoption of optimized rotations for precise quantization of LLMs.

4.3.2 Compatibility with GPTQ

In the context where both weights and activations are quantized, we observed that the learned rotations tend to adapt effectively to both weight and activation quantization. Given that GPTQ significantly helps mitigate the errors due to weight quantization, but leaves activation quantization untouched, we elect to optimize the rotation matrices with respect to a network where only activations are quantized. This approach allows the rotation to more efficiently manage the activation quantization error while leaving the weight quantization error to be addressed by GPTQ. As shown in Table 3, this modification resulted in superior performance in both W4A4 and W4A4KV4 settings in the LLaMA-2 7B model, which is the configuration we have chosen to utilize throughout the rest of this paper.

4.3.3 Rotation type

In Table 4, we evaluate the impact of random orthogonal floating-point rotation matrices and random Hadamard matrices on quantization accuracy, utilizing round-to-nearest quantization for our analysis. Prior to optimization, the Hadamard matrices yield a better-quantized network performance compared to floating-point rotation matrices. However, after optimization, the initial choice of rotation, whether floating-point or Hadamard, becomes less significant. This is likely due to the loss-aware rotation optimization’s ability to locate an optimal local minima that effectively minimizes quantization error, thereby enhancing robustness to varying types of rotation initialization.

4.3.4 Comparison with QuaRot

Compared to QuaRot [5], which exhibits significant accuracy variances in quantized networks—experiencing drops of 28.1 and 33.2 points when quantizing a 70B model with round-to-nearest methods to W4A4 and W4A4KV4—this degradation stems from inherent noise in using random rotation matrix that introduce high variance and compromise robustness. In contrast, SpinQuant_{had} consistently maintains high accuracy across various configurations, achieving improvements of 2.0 to 28.6 points over QuaRot (Table 5), while utilizing fewer online Hadamard matrices (two per block in SpinQuant_{had} versus four per block in QuaRot).

Furthermore, the integration of R_2 in SpinQuant effectively reduces in-block outliers, thereby enabling SpinQuant_{no had} to deliver optimal performance in W4A8 settings. SpinQuant_{no had} can be achieved by simply substituting the model weights with rotated weights, making it a more straightforward and efficient approach compared to QuaRot, which requires modifying the model architecture and special kernel support.

Table 4: Floating-point(FP) rotation vs Hadamard rotation on a LLaMA-2 7B model.

#Bits (W-A-KV)	Task	No Cayley + RTN		Cayley + RTN	
		FP	Hadamard	FP init.	Hadamard init.
4-16-16	0-shot ⁸ Avg.(\uparrow) Wiki(\downarrow)	62.5 \pm 0.8 6.7 \pm 0.12	62.4 \pm 1.0 6.9 \pm 0.45	64.9 \pm 0.4 5.5 \pm 0.01	64.6 \pm 0.3 5.5 \pm 0.01
4-4-16	0-shot ⁸ Avg.(\uparrow) Wiki(\downarrow)	49.4 \pm 2.8 15.9 \pm 4.04	59.0 \pm 1.0 8.2 \pm 0.73	61.6 \pm 0.4 6.2 \pm 0.06	61.8 \pm 0.4 6.1 \pm 0.03
4-4-4	0-shot ⁸ Avg.(\uparrow) Wiki(\downarrow)	48.3 \pm 2.7 18.2 \pm 4.35	58.7 \pm 1.0 8.2 \pm 0.36	61.5 \pm 0.8 6.3 \pm 0.08	61.5 \pm 0.3 6.2 \pm 0.03

Table 5: Comparison with QuaRot [5].

	LLaMA-3 8B (FP: 69.6, 6.1)				LLaMA-3 70B (FP: 74.5, 2.8)			
	4-4-16		4-4-4		4-4-16		4-4-4	
0-shot ⁸ Wiki Avg.(\uparrow) (\downarrow)	59.5 64.6	10.4 7.7	58.6 64.1	10.9 7.8	41.5 70.1	91.2 4.1	41.3 70.1	92.4 4.1
QuaRot+RTN	-	-	-	-	-	-	-	-
SpinQuant _{had} +RTN	-	-	-	-	-	-	-	-
QuaRot+GPTQ	63.8	7.9	63.3	8.0	65.4	20.4	65.1	20.2
SpinQuant _{had} +GPTQ	65.8	7.1	65.5	7.3	69.5	5.5	69.3	5.5

4.4 Illustrative analysis of the rotation efficacy

The rationale behind rotating network weights and activations can be elucidated through a straightforward example. Consider an activation (X) represented as a 2D vector, where one entry x_1 consistently receives higher magnitude activations than x_2 (as depicted in Figure 6(a)). Quantizing these components together typically results in a quantization range dominated by x_1 , thereby compromising the precision for x_2 .

From an information entropy standpoint, expanding each axis to fully utilize the available quantization range maximizes the representational capacity of each axis. Thus, matrix rotation emerges as an intuitive solution. In a 2D scenario, rotating the axis by 45° equalizes the value representation range across axes (illustrated in Figure 6(b)). Assuming the network as a black box without knowledge of the exact activation distribution, uniformly rotating all axes by the maximal degree (45° in 2D) can optimize distribution evenness across each axis, partially explaining why Hadamard rotation often outperforms random rotation matrices.

Taking this further, if the activation distribution is known, treating the network as a white box during quantization allows for the identification of more optimal rotations than Hadamard. For instance, in a 3D scenario depicted in Figure 6(c-d), where x_1 's magnitude is four times that of x_2 and x_3 , rotating the distribution by 45° along x_3 and x_2 redistributes the maximum values from [2, 0.5, 0.5] to [1, 1, 1.414]. However, even more optimal rotation strategies may exist, and learning the rotation can help pinpoint the most effective rotation for a given distribution.

This opens up intriguing research avenues, such as determining if, given an activation distribution with known outlier axes and magnitudes, a closed-form solution for the optimal rotation matrix that

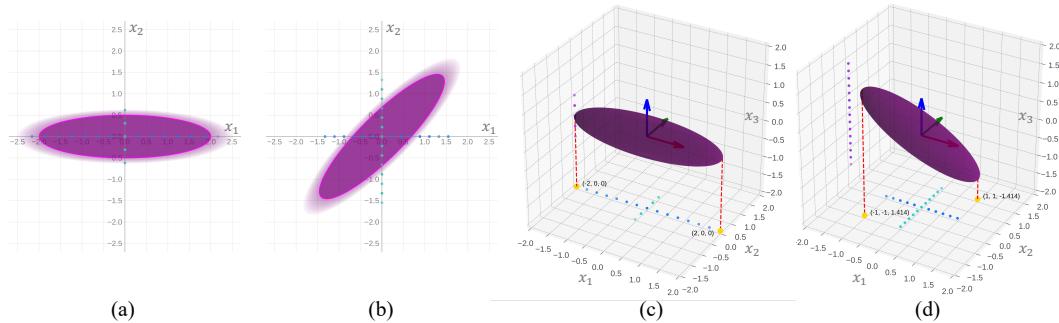


Figure 6: An illustration of how rotation helps reduce outliers and maximize quantization range utilization.

Table 6: Real-time end-to-end speed measurement of LLaMA-3 8B on MacBook M1 Pro CPU.

Method	#Bits _(W-A)	Decoding speed
FloatingPoint	16-16	177.15 ms/token
SpinQuant _{no had}	4-8	58.88 ms/token
SpinQuant _{had}	4-8	63.90 ms/token

evenly distributes magnitude across different axes can be derived. Additionally, it raises the question of whether this theoretically calculated rotation yields the best quantization performance. We leave this question to future research.

4.5 Speed measurement

We conduct an end-to-end speed measurement of the LLaMA-3 8B model with W16A16 and W4A8 configurations on a MacBook M1 Pro CPU (OS version 14.5). The results in Table 6 demonstrate that 4-bit quantization yields a $\sim 3 \times$ speedup compared to the 16-bit model. Comparing SpinQuant_{had} to SpinQuant_{no had}, online Hadamard processing introduced a modest 8% increase in latency. Therefore, it is a trade-off between using SpinQuant_{no had} without online Hadamard for its simpleness or using SpinQuant_{had} with online Hadamard rotations for higher accuracy in lower-bit activation quantization. Detailed GPU latency results are provided in the Appendix.

5 Related Work

Quantization Neural network quantization has been demonstrated as an effective tool for model size compression and storage reduction [31, 19, 30, 22]. However, in large language models (LLMs), quantization presents unique challenges due to the presence of numerous outliers. These outliers dominate the quantization range, leaving only a few effective bits for the majority of values. Various strategies have been proposed to address the difficulties in LLM quantization. These include separating outliers and using mixed precision [11, 43, 18, 15, 12], employing Hessian-based methods to mitigate quantization difficulty [14], trading outliers between weights and activations [46, 23, 25] utilizing weight equalization [30], outlier suppression [44, 45], channel reassembly [24] and even suggesting architectural modifications to handle outliers during pre-training[48]. Recently two QuIP papers [7, 42] introduce the incoherence processing using random rotation matrices and applying vector quantization on the weights for compression. This does introduce extra overhead and imposes some constraints on the devices the LLM is deployed to in the availability of vector quantization kernels.

Optimization in orthonormal space The optimization of rotation matrices is carried out within the *Stiefel Manifold* [16], which encompasses all orthonormal matrices. Optimization while staying on this manifold can be done by *e.g.*, parameterizing a skew-symmetric matrix and applying the *Cayley* transformation on top of it [32], or using a matrix exponential [1, 20]. However, these methods rely on expensive inverse or matrix-exponential functions that are applied every iteration. Instead, we follow the more efficient method named *Cayley SGD* [21], which can be applied to optimize a rotation matrix R for arbitrary loss functions efficiently. *Cayley SGD* relies on an iterative approximation of the *Cayley* Transform that is conducted solely with matrix multiplications.

6 Conclusions

In this paper, we present SpinQuant, a novel quantization technique that utilizes learned rotation to effectively bridge the performance gap between full precision and 4-bit weight, activation, and kv-cache quantization. At its core, SpinQuant leverages the rotation invariance property of LLM models to insert rotation matrices that diminish outliers in the weights and intermediate activations while maintaining the network’s full-precision output numerically identical. Additionally, SpinQuant incorporates *Cayley SGD* for optimizing rotation matrices, resulting in improved and robust quantization outcomes. Importantly, SpinQuant is compatible with more advanced weight quantization techniques (*e.g.*, GPTQ) and demonstrates state-of-the-art performance.

7 Acknowledgement

We extend our gratitude to Scott Wolchok and Chen Lai for their crucial contributions to latency measurement on the MacBook M1 Pro CPU, and to Geonhwa Jeong and Jiecao Yu for their expert support in GPU latency assessment.

References

- [1] P-A Absil and Jérôme Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on Optimization*, 22(1):135–158, 2012.
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- [3] AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md.
- [4] Saleh Ashkboos, Maximilian L Croci, Marcelo Gennari do Nascimento, Torsten Hoefer, and James Hensman. Slicept: Compress large language models by deleting rows and columns. In *The Twelfth International Conference on Learning Representations*, 2023.
- [5] Saleh Ashkboos, Amirkeivan Mohtashami, Maximilian L. Croci, Bo Li, Martin Jaggi, Dan Alistarh, Torsten Hoefer, and James Hensman. Quarot: Outlier-free 4-bit inference in rotated llms. 2023.
- [6] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical common-sense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [7] Jerry Chee, Yaohui Cai, Volodymyr Kuleshov, and Christopher M De Sa. Quip: 2-bit quantization of large language models with guarantees. *Advances in Neural Information Processing Systems*, 36, 2024.
- [8] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. *arXiv preprint arXiv:1905.10044*, 2019.
- [9] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.
- [10] Samuel Rhys Cox and Wei Tsang Ooi. Conversational interactions with npcs in llm-driven gaming: Guidelines from a content analysis of player feedback. In *International Workshop on Chatbot Research and Design*, pages 167–184. Springer, 2023.
- [11] Tim Dettmers, Mike Lewis, Younes Belkada, and Luke Zettlemoyer. Llm.int8(): 8-bit matrix multiplication for transformers at scale. 2022.
- [12] Vage Egiazarian, Andrei Panferov, Denis Kuznedelev, Elias Frantar, Artem Babenko, and Dan Alistarh. Extreme compression of large language models via additive quantization. *arXiv preprint arXiv:2401.06118*, 2024.
- [13] Nelson Elhage, Robert Lasenby, and Christopher Olah. Privileged bases in the transformer residual stream. *Transformer Circuits Thread*, 2023.
- [14] Elias Frantar, Saleh Ashkboos, Torsten Hoefer, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers. *arXiv preprint arXiv:2210.17323*, 2022.
- [15] Wei Huang, Haotong Qin, Yangdong Liu, Yawei Li, Xianglong Liu, Luca Benini, Michele Magno, and Xiaojuan Qi. Slim-llm: Salience-driven mixed-precision quantization for large language models. *arXiv preprint arXiv:2405.14917*, 2024.
- [16] Ioan Mackenzie James. *The topology of Stiefel manifolds*, volume 24. Cambridge University Press, 1976.
- [17] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

- [18] Sehoon Kim, Coleman Hooper, Amir Gholami, Zhen Dong, Xiuyu Li, Sheng Shen, Michael W Mahoney, and Kurt Keutzer. Squeezellm: Dense-and-sparse quantization. *arXiv preprint arXiv:2306.07629*, 2023.
- [19] Raghuraman Krishnamoorthi. Quantizing deep convolutional networks for efficient inference: A whitepaper. *arXiv preprint arXiv:1806.08342*, 2018.
- [20] Mario Lezcano-Casado and David Martinez-Rubio. Cheap orthogonal constraints in neural networks: A simple parametrization of the orthogonal and unitary group. 2019.
- [21] Jun Li, Li Fuxin, and Sinisa Todorovic. Efficient riemannian optimization on the stiefel manifold via the cayley transform. *arXiv preprint arXiv:2002.01113*, 2020.
- [22] Yuhang Li, Ruihao Gong, Xu Tan, Yang Yang, Peng Hu, Qi Zhang, Fengwei Yu, Wei Wang, and Shi Gu. Brecq: Pushing the limit of post-training quantization by block reconstruction. In *International Conference on Learning Representations (ICLR)*, 2021.
- [23] Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Xingyu Dang, and Song Han. Awq: Activation-aware weight quantization for llm compression and acceleration. *arXiv preprint arXiv:2306.00978*, 2023.
- [24] Jing Liu, Ruihao Gong, Xiuying Wei, Zhiwei Dong, Jianfei Cai, and Bohan Zhuang. Qllm: Accurate and efficient low-bitwidth quantization for large language models. *arXiv preprint arXiv:2310.08041*, 2023.
- [25] Shih-yang Liu, Zechun Liu, Xijie Huang, Pingcheng Dong, and Kwang-Ting Cheng. Llm-fp4: 4-bit floating-point quantized transformers. *arXiv preprint arXiv:2310.16836*, 2023.
- [26] Zechun Liu, Barlas Oguz, Changsheng Zhao, Ernie Chang, Pierre Stock, Yashar Mehdad, Yangyang Shi, Raghuraman Krishnamoorthi, and Vikas Chandra. Llm-qat: Data-free quantization aware training for large language models. *arXiv preprint arXiv:2305.17888*, 2023.
- [27] Zechun Liu, Changsheng Zhao, Forrest Iandola, Chen Lai, Yuandong Tian, Igor Fedorov, Yunyang Xiong, Ernie Chang, Yangyang Shi, Raghuraman Krishnamoorthi, et al. Mobilellm: Optimizing sub-billion parameter language models for on-device use cases. *arXiv preprint arXiv:2402.14905*, 2024.
- [28] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.
- [29] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. *arXiv preprint arXiv:1809.02789*, 2018.
- [30] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1325–1334, 2019.
- [31] Markus Nagel, Rana Ali Amjad, Mart Van Baalen, Christos Louizos, and Tijmen Blankevoort. Up or down? Adaptive rounding for post-training quantization. In *International Conference on Machine Learning (ICML)*, 2020.
- [32] Yasunori Nishimori and Shotaro Akaho. Learning algorithms utilizing quasi-geodesic flows on the stiefel manifold. *Neurocomputing*, 67:106–135, 2005.
- [33] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [34] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémie Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
- [35] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- [36] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan LeBras, and Yejin Choi. Socialqa: Commonsense reasoning about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- [37] Wenqi Shao, Mengzhao Chen, Zhaoyang Zhang, Peng Xu, Lirui Zhao, Zhiqian Li, Kaipeng Zhang, Peng Gao, Yu Qiao, and Ping Luo. Omnipquant: Omnidirectionally calibrated quantization for large language models. *arXiv preprint arXiv:2308.13137*, 2023.

- [38] Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- [39] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabilan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. Large language models in medicine. *Nature medicine*, 29(8):1930–1940, 2023.
- [40] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [41] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- [42] Albert Tseng, Jerry Chee, Qingyao Sun, Volodymyr Kuleshov, and Christopher De Sa. Quip#: Even better lilm quantization with hadamard incoherence and lattice codebooks. *arXiv preprint arXiv:2402.04396*, 2024.
- [43] Mart van Baalen, Markus Nagel Andrey Kuzmin, Peter Couperus, Cedric Bastoul, Eric Mahurin, Tijmen Blankevoort, and Paul Whatmough. Gptvq: The blessing of dimensionality for lilm quantization. 2023.
- [44] Xiuying Wei, Yunchen Zhang, Xiangguo Zhang, Ruihao Gong, Shanghang Zhang, Qi Zhang, Fengwei Yu, and Xianglong Liu. Outlier suppression: Pushing the limit of low-bit transformer language models. *arXiv preprint arXiv:2209.13325*, 2022.
- [45] Xiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and Xianglong Liu. Outlier suppression+: Accurate quantization of large language models by equivalent and optimal shifting and scaling. *arXiv preprint arXiv:2304.09145*, 2023.
- [46] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *CVPR*, 2022.
- [47] Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. Smoothquant: Accurate and efficient post-training quantization for large language models. In *International Conference on Machine Learning*, pages 38087–38099. PMLR, 2023.
- [48] Tijmen Blankevoort Yelysei Bondarenko, Markus Nagel. Quantizable transformers: Removing outliers by helping attention heads do nothing. 2023.
- [49] Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*, 2019.
- [50] Yilong Zhao, Chien-Yu Lin, Kan Zhu, Zihao Ye, Lequn Chen, Size Zheng, Luis Ceze, Arvind Krishnamurthy, Tianqi Chen, and Baris Kasikci. Atom: Low-bit quantization for efficient and accurate lilm serving. *arXiv preprint arXiv:2310.19102*, 2023.

A Appendix / supplemental material

A.1 Complete results of main result table

In Tables 7, 8 and 9, we show the complete results of Table 1. We compare the accuracy on eight zero-shot commonsense reasoning tasks including ARC-easy, ARC-challenge [9], BoolQ [8], PIQA [6], SIQA [36], HellaSwag [49], OBQA [29], and WinoGrande [35] as well as the perplexity score on WikiText2 testset [28]. We compare our results with previous works including SmoothQuant[46], LLM-QAT[26], GPTQ [14], OmniQuant [37], QuIP# [42].

A.2 Results on 3-bit weight quantization

We present the 3-bit weight and 8-bit activation quantization results across seven models in Table 10. Our method, SpinQuant, successfully reduces the gap to the full-precision network from the previous 9.0 – 28.0 points to 1.2 – 5.3 points, demonstrating its effectiveness for low-bit quantization.

A.3 Cayley optimization choice

In Table 11, we evaluate the impact of varying the number of samples and iterations used in *Cayley* optimization. Given the limited trainable parameters in the rotation matrix and its constraint optimization nature, minimal calibration data and iterations are sufficient to optimize the rotation for better quantization. The findings indicate that rotation optimization is resilient to modifications in the number of samples. Even though we used 800 samples in our experiments, reducing this to 128 samples does not lead to a significant change in the perplexity. Furthermore, we examined the optimal number of iterations and found that the wiki perplexity ceases to decrease and stabilizes at 100 iterations. Consequently, we chose to use 100 iterations in all our experiments.

A.4 Quantization choice

We conduct an ablation study on symmetric vs asymmetric quantization and whether to clip the min-max ranges or not during activation and KV-cache quantization. The results in Table 12 show that for both activation quantization and KV-cache quantization, asymmetric quantization outperforms symmetric quantization. In the clip settings, we set the activation clipping ratio to 0.9 and the KV-cache clipping ratio to 0.95 as suggested in the previous works [50]. However, the results show that clipping the range or not does not impact the final result significantly. Therefore we opt for no clipping, *i.e.*, using the min-max quantization for activation and KV cache quantization across our experiments due to its simplicity.

A.5 Calibration data choice

To assess the robustness of SpinQuant with respect to calibration data used in rotation optimization we use C4 dataset [33] as calibration data and performe experiments on the LLaMA-2 7B model. The results in Table 13 reflect that using C4 datasets yields consistent results with utilizing the Wiki dataset, showing that SpinQuant is robust to calibration data choice.

A.6 Latency measurement on GPU

We measured the latency of each component in LLaMA-3 70B decoding with weight quantization to FP8 using SpinQuant_{had} . As shown in Figure 7, online Hadamard rotation accounts for only 3.6% of total computation latency.

B Analysis

B.1 Gradient Analysis

On the one hand, we have shown that the class of LLMs we are interested in are rotation invariant, *i.e.* the full-precision model output does not change regardless of what R is. On the other hand, we are claiming that some R are better than others for quantized LLM and that better R can be learned with

Table 12: Ablation of symmetric and asymmetric quantization and range clipping options on LLaMA-2 7B.

#Bits (W-A-KV)	K asym	K clip	A asym	A clip	RTN	GPTQ
4-4-16	–	–	✗	✗	61.2 ±0.6	6.3
4-4-16	–	–	✓	✗	61.8 ±0.4	6.1
4-4-16	–	–	✓	✓	62.1 ±0.6	6.0
4-4-4	✗	✗	✓	✗	61.4 ±0.5	6.2
4-4-4	✓	✗	✓	✗	61.5 ±0.6	6.2
4-4-4	✓	✓	✓	✗	61.5 ±0.3	6.2

Table 13: Ablation study on calibration data choice using LLaMA-2 7B.

Calibration Data	#Bits (W-A-KV)	ARC-e (↑)	ARC-c (↑)	BoolQ (↑)	PIQA (↑)	SIQA (↑)	HellaS. (↑)	OBQA (↑)	WinoG. (↑)	Avg. (↑)	Wiki2 (↓)
Wiki2	4-4-16	72.1	47.5	74.4	77.0	47.3	73.2	54.4	66.9	64.1	5.9
Wiki2	4-4-4	72.6	47.5	73.9	77.0	47.2	73.0	54.1	66.9	64.0	5.9
C4	4-4-16	72.5	47.3	74.8	77.6	47.7	73.2	55.4	66.2	64.3	5.9
C4	4-4-4	72.5	47.9	74	78.4	46.7	73.1	55.5	66.4	64.3	6

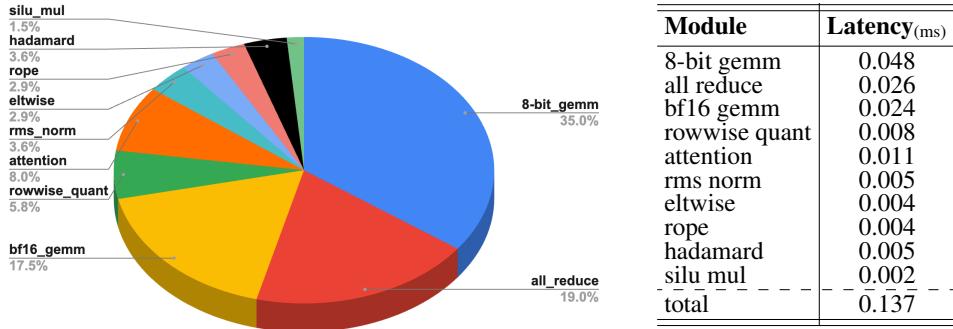


Figure 7: The latency measurement of LLaMA-3 70B model with 8 H100 GPUs

backpropagation on (2). To reconcile these seemingly conflicting claims, we inspect the gradient of the output of a single linear, W , and activations, X , which are both rotated and quantized:

$$\frac{\partial \sum_{ij} (Q(WR^{-1})Q(RX))_{ij}}{\partial R_{mn}} = \sum_{ij} -(WR^{-1})_{im}(R^{-1}Q(RX))_{nj} + Q(WR^{-1})_{im}X_{nj} \quad (5)$$

We see that equation (5):

- is non-zero in general, which validates our approach of using backpropagation to learn R
- reduces to 0 when quantization is not present, which validates the claim that it only makes sense to learn R for quantized models
- demonstrates that two components move the gradient with respect to R away from 0: 1) differences in quantized and unquantized rotated weights; 2) differences in quantized and unquantized rotated activations

B.2 Loss Analysis

While Sec. 4 shows that learning R yields significant benefits on zero-shot reasoning tasks, in this section we shed some light on why our method is able to achieve accuracy gains. Intuitively, we expect the end-to-end signal to (quantization) noise ratio (SNR) to improve as a result of learning R . In other words, learning R should bring the quantized model output closer to the floating point model output. As Table 14 shows, we observe an SNR improvement of 3.8 dB when introducing a random R into LLaMA-2 7B with weights/activations quantized to 4 bits, and then an additional 5.9dB improvement after learning R , all measured on the WikiText2 [28] test set. Figure 8a shows that the batch-level training set SNR during R training progressively improves as expected, as well as the layer-level SNR for a particular layer in Figure 8b. Digging a bit deeper, Figure 8c shows the layer-level SNR improvement for each layer as a result of training R . We see that, perhaps

Table 14: Average end-to-end signal to quantization noise ratio (dB) for LLaMA-2 7B with weights and activations quantized to 4 bits on wiki2 test set

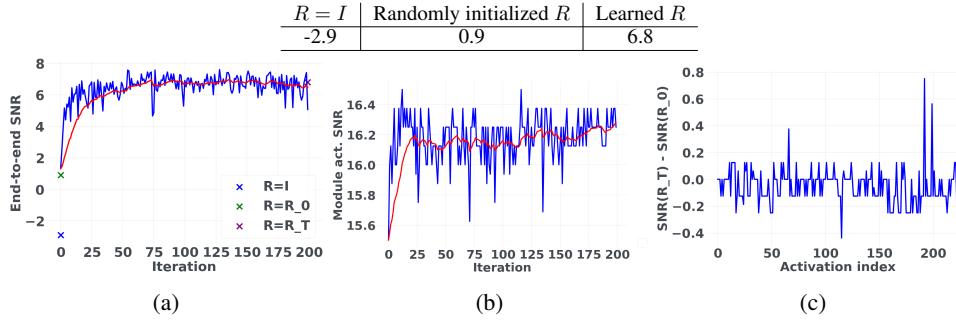


Figure 8: Training curves for LLaMA-2 7B with 4-bit weights and 4-bit activations in wiki2 train set. (a) End-to-end quantization SNR. R_0 and R_T denote randomly initialized rotation and learned rotation after $T = 200$ iterations; (b) Activation quantization. SNR for layer 27 attention out projection; (c) Improvement in activation quantization SNR after optimization of R for each layer.

counter-intuitively, layer-level SNR improves significantly for a few layers, but does not change much for most layers, and even gets worse for one of the layers. We hypothesize that: 1) certain layers have a disproportionate impact on model output or have a disproportionately low quantization SNR without rotation; 2) The process of optimizing R rotates the residual stream basis such as to prioritize improving the SNR of such layers, possibly at the cost of hurting less important layers.

C Distribution visualizations before and after rotation

We present visualizations of the activation distributions before and after rotation in Figures 9 and 10, respectively. Similarly, the weight distributions before and after rotation are depicted in Figures 11 and 12. Overall, after rotation, the extreme values are attenuated, and the distribution exhibits no noteworthy outliers across the token dimension. Additionally, we make an interesting observation: in several activation layers, the first token displays substantial values in multiple channels. After rotation, this outlier is distributed across all channels of the first token. Although per-token activation quantization can readily manage this distribution, investigating the source of these outliers and reducing them prior to applying SpinQuant might further enhance quantization accuracy, which could be a potential future research direction.

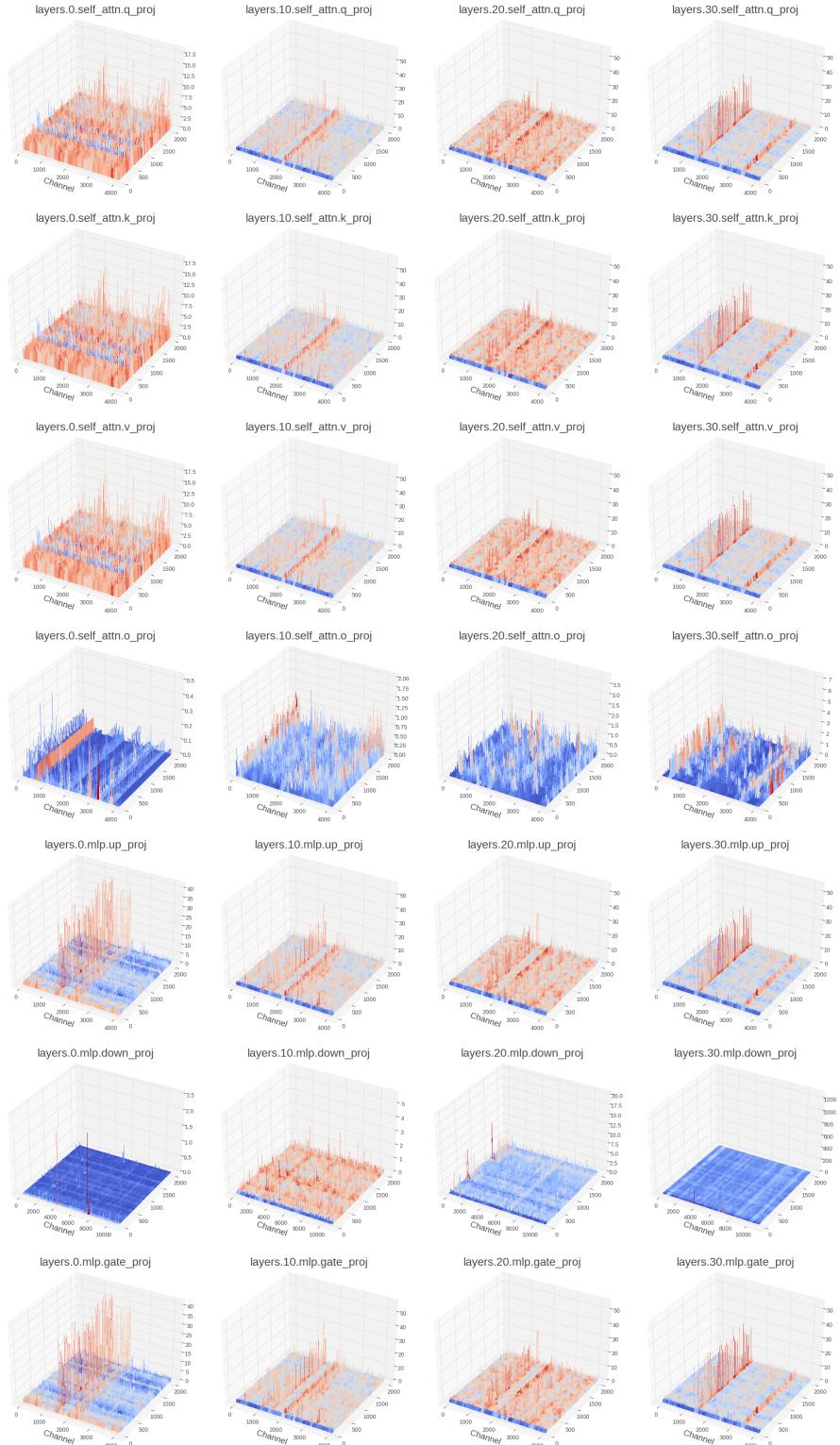


Figure 9: Magnitude of the input activations of a linear layer in {1st, 11th, 21st, and 31st} blocks in LLaMA-2 7B model **before rotation**.

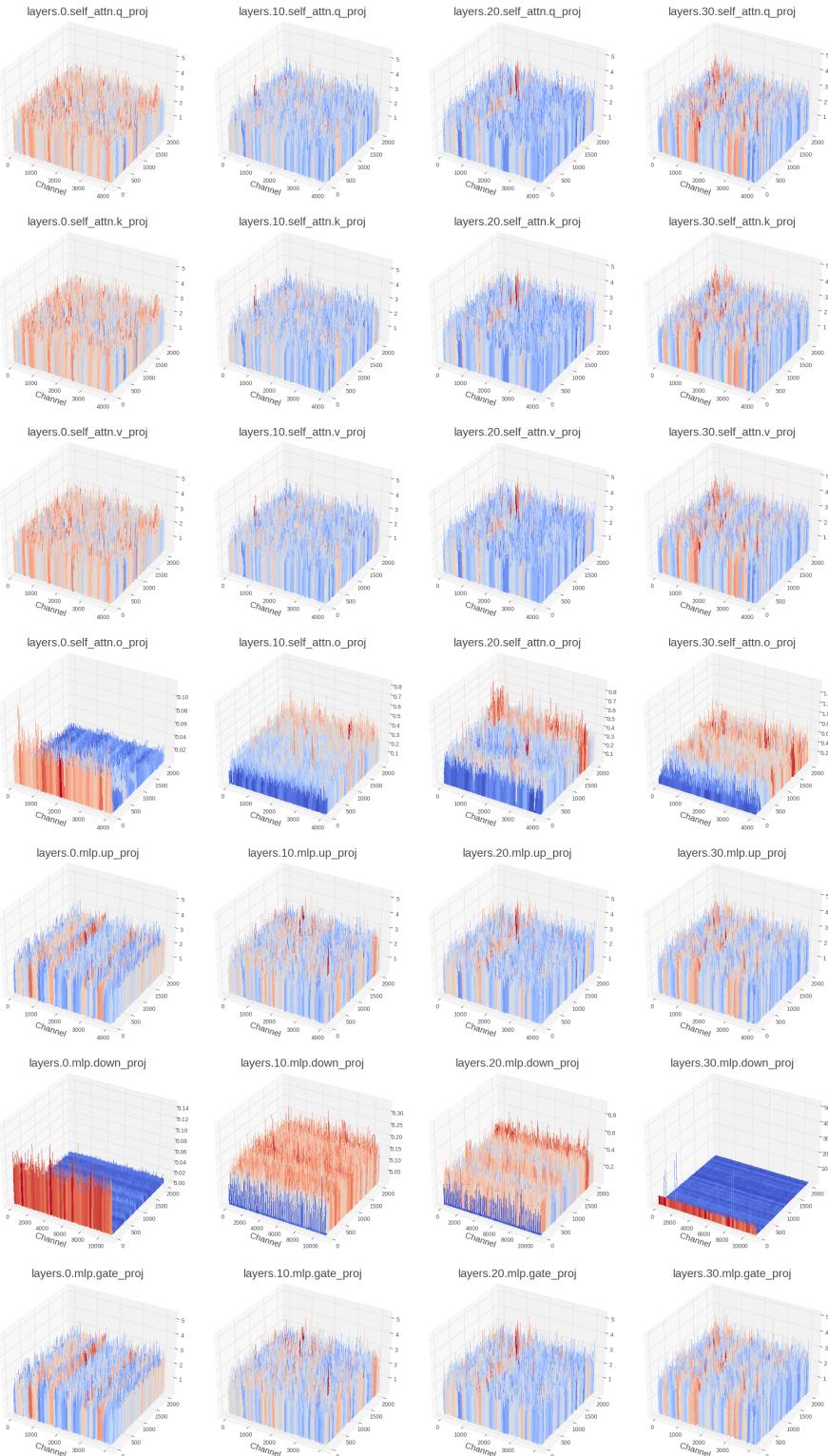


Figure 10: Magnitude of the input activations of a linear layer in {1st, 11th, 21st, and 31st} blocks in LLaMA-2B model **after rotation**.

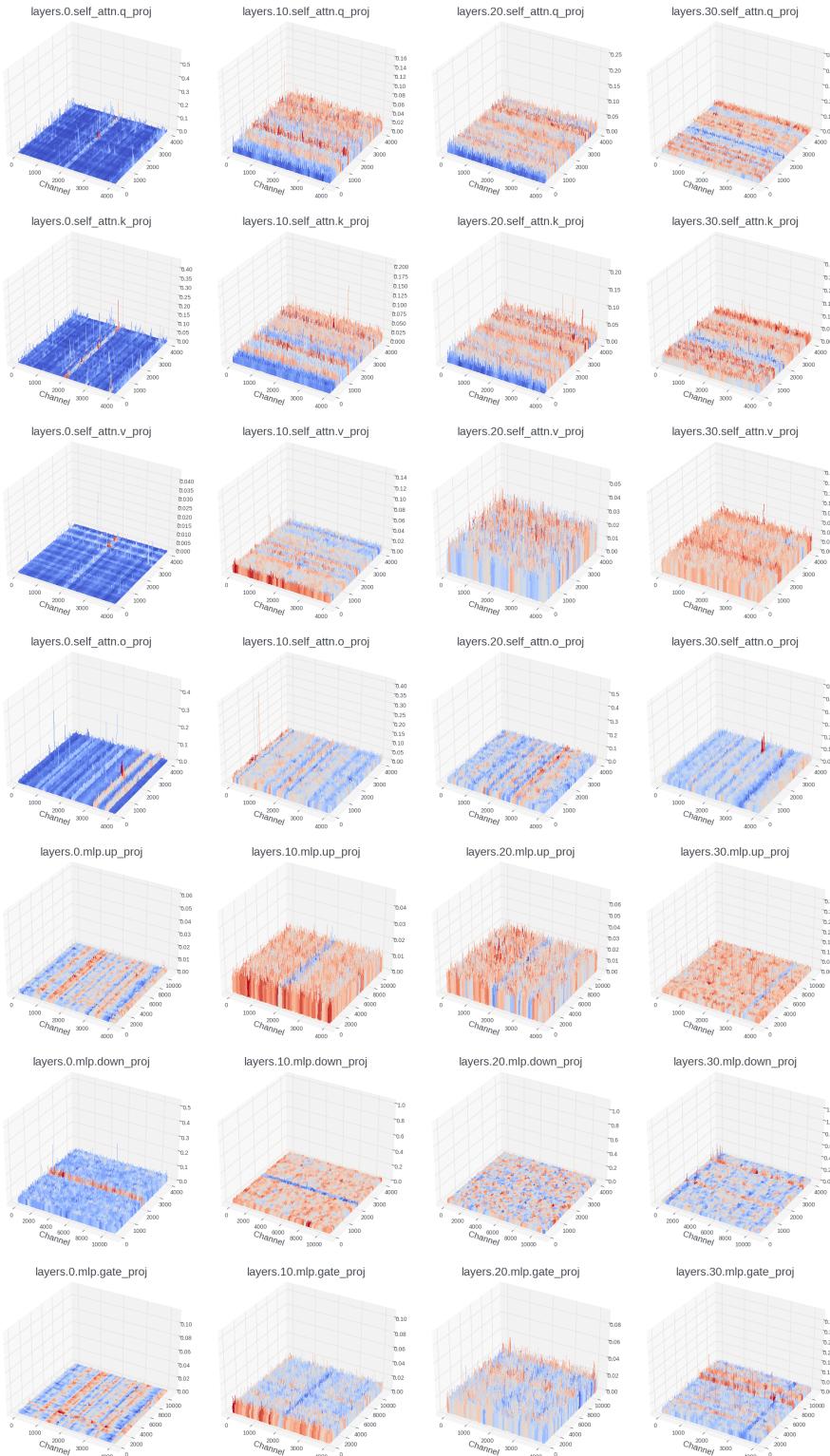


Figure 11: Magnitude of the weights of a linear layer in {1st, 11th, 21st, and 31st} blocks in LLaMA-2 7B before rotation.

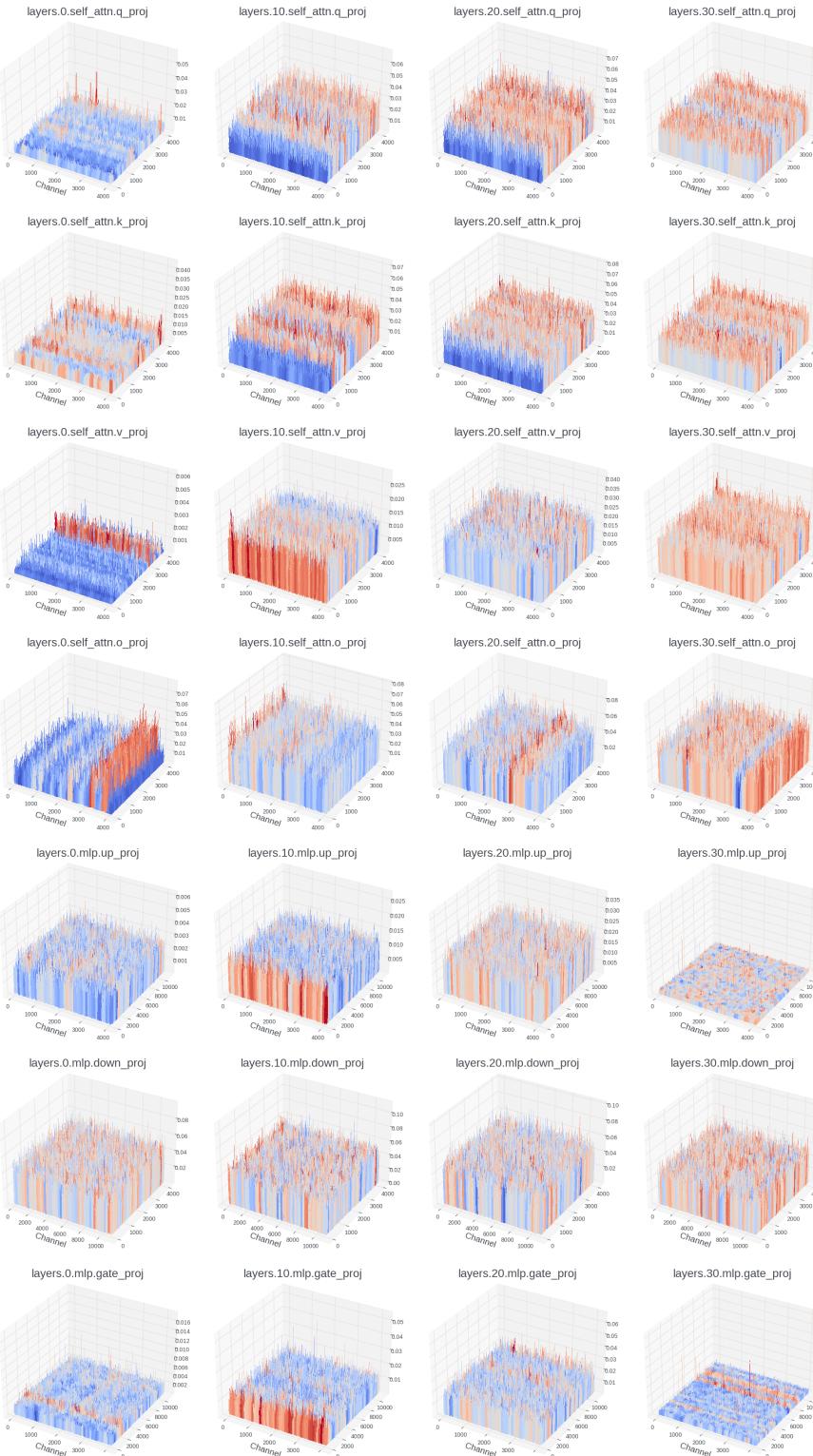


Figure 12: Magnitude of the weights of a linear layer in $\{1^{st}, 11^{th}, 21^{st}, \text{ and } 31^{st}\}$ blocks in LLaMA-2 7B after rotation.